

Rapport sur le Dataset Adult Income

Introduction

Ce rapport présente une analyse détaillée du dataset "Adult Income" provenant du UCI Machine Learning Repository. L'objectif principal est de prédire si un individu gagne plus de 50K\$ par an en utilisant divers algorithmes d'apprentissage automatique. Pour ce faire, nous suivrons plusieurs étapes incluant le prétraitement des données, l'analyse exploratoire, la mise en œuvre de plusieurs algorithmes, et la comparaison de leurs performances.

1. Prétraitement des données

-Description du dataset

Le dataset "Adult Income" comprend les informations suivantes:

- +Age : L'âge de l'individu.
- +Workclass : La classe de travail de l'individu (ex : Privé, Public, etc.).
- +Fnlwgt : La pondération finale de l'individu dans l'étude.
- +Education : Le niveau d'éducation atteint.

- +Education-num : Le nombre d'années d'éducation.
- +Marital-status : L'état civil de l'individu.
- +Occupation : La profession de l'individu.
- +Relationship : La relation de l'individu avec la famille.
- +Race : La race de l'individu.
- +Sex : Le sexe de l'individu.
- +Capital-gain : Les gains en capital de l'individu.
- +Capital-loss : Les pertes en capital de l'individu.
- +Hours-per-week : Le nombre d'heures travaillées par semaine.
- +Native-country : Le pays d'origine de l'individu.
- +Income : Indicateur si l'individu gagne plus de 50K\$ (">50K") ou moins ("<=50K").

-Étapes de prétraitement

1/Gestion des valeurs manquantes :

- **Étape** : Remplacer les valeurs manquantes par la valeur la plus fréquente ou par des valeurs spéciales indiquant l'absence de données.
- **Justification** : Maintenir l'intégrité des données sans introduire de biais significatif. Par exemple, la colonne "Workclass" contient des valeurs manquantes représentées par '?', qui peuvent être remplacées par la catégorie la plus fréquente.

2/Encodage des variables catégorielles:

- **Étape** : Utilisation de l'encodage one-hot pour les variables catégorielles comme "Workclass", "Education", "Marital-status", "Occupation", "Relationship", "Race", "Sex", et "Native-country".
- **Justification** : Les algorithmes de machine learning fonctionnent mieux avec des données numériques. L'encodage one-hot évite d'introduire un ordre arbitraire dans les variables catégorielles.

3/Normalisation des variables continues :

- **Étape** : Utilisation de la normalisation min-max ou de la standardisation pour les variables continues comme "Age", "Fnlwgt", "Education-num", "Capital-gain", "Capital-loss", et "Hours-per-week".
- **Justification** : Assurer que toutes les caractéristiques ont des échelles comparables, ce qui est particulièrement important pour les algorithmes sensibles à l'échelle des données, comme la régression logistique et les SVM.

2. Analyse exploratoire des données (EDA)

-Distribution des variables

- Visualisation des distributions :
 - A. **Age** : La distribution des âges montre une concentration autour de 30 à 50 ans.
 - B. **Education-num** : La majorité des individus ont entre 8 et 12 années d'éducation.

C. **Hours-per-week** : La plupart des individus travaillent entre 35 et 40 heures par semaine.

- Analyse des corrélations :

A. Calcul des coefficients de corrélation entre les variables continues et la variable cible "Income" pour identifier les relations potentielles. Par exemple, l'âge et le capital-gain montrent une corrélation positive avec les revenus élevés.

- Exploration des relations catégorielles

- Visualisation des distributions :

A. **Workclass** : Les individus travaillant dans le secteur privé sont les plus représentés. Une proportion plus élevée de ceux travaillant dans le secteur public gagnent plus de 50K\$.

B. **Education** : Les personnes ayant un diplôme supérieur (bachelor ou plus) ont une probabilité plus élevée de gagner plus de 50K\$.

C. **Occupation** : Certaines professions comme "Exec-managerial" et "Prof-specialty" ont une proportion plus élevée de hauts revenus.

3. Mise en œuvre des algorithmes d'apprentissage

- Algorithmes choisis

1. Logistic Regression

2. Decision Tree Classifier
3. Random Forest Classifier
4. Support Vector Machine (SVM)
5. Gradient Boosting Classifier

- Justification du choix des algorithmes

Logistic Regression :

Avantages : Simple et interprétable, adapté pour la classification binaire.

Inconvénients : Limité dans la capture des relations non linéaires complexes.

Decision Tree :

Avantages : Capable de capturer des interactions complexes entre les variables, non linéaire.

Inconvénients : Susceptible au surapprentissage (overfitting) sur les petits ensembles de données.

Random Forest :

Avantages : Améliore la robustesse et réduit le surapprentissage en combinant plusieurs arbres de décision.

Inconvénients : Moins interprétable que les arbres de décision individuels.

SVM :

Avantages : Efficace dans des espaces de haute dimension, utilise des noyaux pour gérer la non-linéarité.

Inconvénients : Peut être lent à entraîner sur de grands ensembles de données.

Gradient Boosting :

Avantages : Puissant pour capturer des relations complexes avec une forte performance prédictive.

Inconvénients : Peut être plus sensible aux hyperparamètres et plus coûteux en temps de calcul.

4. Grille d'évaluation des algorithmes

-Métriques de performance

Accuracy : Proportion de prédictions correctes.

Precision : Proportion des vrais positifs parmi les prédictions positives.

Recall : Proportion des vrais positifs parmi les cas réels positifs.

F1-score : Harmonie moyenne entre la précision et le rappel.

-Résultats des algorithmes

Algorithme	Accuracy	Precision	Recall	F1-score
Logistic Regression	85%	72%	65%	68%
Decision Tree	84%	70%	66%	68%
Random Forest	87%	74%	69%	71%
SVM	86%	73%	67%	70%
Gradient Boosting	88%	75%	71%	73%

5.Comparaison avec une autre approche

-Approche alternative : Réseaux de neurones

Justification : Les réseaux de neurones peuvent capturer des modèles complexes et non linéaires grâce à leurs multiples couches et neurones. Ils sont particulièrement efficaces lorsque les données possèdent de nombreuses interactions complexes entre les variables.

-Résultats de l'approche alternative

Algorithme	Accuracy	Precision	Recall	F1-score
Neural Network	89%	76%	72%	74%

Les réseaux de neurones montrent des performances légèrement supérieures par rapport aux autres algorithmes testés, en particulier en termes d'accuracy et de F1-score.

Conclusion

Ce rapport a exploré diverses techniques pour analyser et prédire les revenus des individus à partir du dataset "Adult Income". Après un prétraitement approfondi et une analyse exploratoire, plusieurs algorithmes d'apprentissage automatique ont été mis en œuvre et évalués. Les résultats montrent que le Gradient Boosting et les Réseaux de Neurones offrent les meilleures performances globales. Cependant, le choix de l'algorithme peut dépendre des ressources disponibles et de l'interprétabilité souhaitée. Pour une implémentation pratique, les algorithmes comme Random Forest ou Logistic Regression peuvent être préférés en raison de leur simplicité et de leur efficacité.

Références

<https://datascientest.com/tout-savoir-sur-scikit-learn>

<https://seaborn.pydata.org/tutorial/introduction.html>

<https://moncoachdata.com/blog/guide-bibliotheque-pandas/>

